

## RNA folding and combinatorial landscapes

Walter Fontana

*Santa Fe Institute, 1660 Old Pecos Trail, Santa Fe, New Mexico 87501  
and Theoretical Division T-13, Los Alamos National Laboratory, Los Alamos, New Mexico 87545*

Peter F. Stadler

*Institut für Theoretische Chemie, Universität Wien, Währingerstraße 17, A-1090 Wien, Austria;  
Max-Planck-Institut für Biophysikalische Chemie, Am Fassberg, D-3400 Göttingen, Germany;  
and Santa Fe Institute, 1660 Old Pecos Trail, Santa Fe, New Mexico 87501*

Erich G. Bornberg-Bauer, Thomas Griesmacher, Ivo L. Hofacker, and Manfred Tacker

*Institut für Theoretische Chemie, Universität Wien, Währingerstraße 17, A-1090 Wien, Austria*

Pedro Tarazona

*Institut für Theoretische Chemie, Universität Wien, Währingerstraße 17, A-1090 Wien, Austria  
and Departamento de Física de la Materia Condensada C-XII, Universidad Autónoma de Madrid, E-28049 Madrid, Spain*

Edward D. Weinberger

*Max-Planck-Institut für Biophysikalische Chemie, Am Fassberg, D-3400 Göttingen, Germany*

Peter Schuster\*

*Institut für Theoretische Chemie, Universität Wien, Währingerstraße 17, A-1090 Wien, Austria;  
Santa Fe Institute, 1660 Old Pecos Trail, Santa Fe, New Mexico 87501;  
and Institut für Molekulare Biotechnologie, Beutenbergstraße 11, D0-6900 Jena, Germany*

(Received 2 October 1992)

In this paper we view the folding of polynucleotide (RNA) sequences as a map that assigns to each sequence a minimum-free-energy pattern of base pairings, known as secondary structure. Considering only the free energy leads to an energy landscape over the sequence space. Taking into account structure generates a less visualizable nonscalar "landscape," where a sequence space is mapped into a space of discrete "shapes." We investigate the statistical features of both types of landscapes by computing autocorrelation functions, as well as distributions of energy and structure distances, as a function of distance in sequence space. RNA folding is characterized by very short structure correlation lengths compared to the diameter of the sequence space. The correlation lengths depend strongly on the size and the pairing rules of the underlying nucleotide alphabet. Our data suggest that almost every minimum-free-energy structure is found within a small neighborhood of any random sequence. The interest in such landscapes results from the fact that they govern natural and artificial processes of optimization by mutation and selection. Simple statistical model landscapes, like Kauffman and Levin's  $n$ - $k$  model [J. Theor. Biol. **128**, 11 (1987)], are often used as a proxy for understanding realistic landscapes, like those induced by RNA folding. We make a detailed comparison between the energy landscapes derived from RNA folding and those obtained from the  $n$ - $k$  model. We derive autocorrelation functions for several variants of the  $n$ - $k$  model, and briefly summarize work on its fine structure. The comparison leads to an estimate for  $k = 7-8$ , independent of  $n$ , where  $n$  is the chain length. While the scaling behaviors agree, the fine structure is considerably different in the two cases. The reason is seen to be the extremely high frequency of neutral neighbors, that is, neighbors with identical energy (or structure), in the RNA case.

PACS number(s): 87.10.+e, 87.15.By, 64.60.Cn

### I. COMBINATORIAL MAPS

Processes like combinatorial optimization and evolutionary adaptation take place on landscapes that result from mapping microconfigurations to energies or nonscalar entities like structures. A classic example from physics is a Hamiltonian that assigns an energy value to a spin configuration. Another instance, taken from operations

research, is a geography that maps tours through a set of cities into transport costs. Many properties of those processes reflect the local and global statistical features of the landscape on which they occur. This leads to the problem of understanding what these features are, and how to study them.

As an example, we study the biologically important landscape induced by the folding of polynucleotides

(RNA). Here, the microconfigurations are sequences over an alphabet of nucleotides, and scalar properties are free energies of secondary or tertiary structure formation. Another scalar could be the rate constant of a reaction involving that structure (e.g., replication of viral RNA). Nonscalar properties like the secondary structure or the 3D (three-dimensional) structure are of particular interest. Here, we consider both the energy as well as the structure landscapes induced by RNA folding, and compare the former with a simple parametrized model landscape, known as the  $n$ - $k$  model [1].

Common to all these examples is a function whose domain is a set of combinatorial complexity—where the elements represent combinations or variations of some kind—and whose range is either  $\mathbb{R}$ , or another set of combinatorial complexity (suitably discretized structures, for example). This motivates the following definition of a combinatory map (CM).

*Definition.* A CM is a quintuple  $(x, d_x; y, d_y; f)$ , where  $x$  and  $y$  are sets endowed with metrics  $d_x$  and  $d_y$ , respectively, and  $f$  is a map  $x \rightarrow y$ . If  $y = \mathbb{R}$  and  $d_y(a, b) = |a - b|$ , we refer to the quintuple as a combinatory landscape (CL).

$(x, d_x)$  is known as configuration space. The natural metric is induced by some physically meaningful set of operations that interconvert configurations. In the present case of RNA's, a configuration (or sequence) space consists of all sequences of fixed length  $n$  over an alphabet of size  $\kappa$  (typically  $\kappa = 4$ ), the interconversion operations are point mutations, and the Hamming metric provides a distance measure between sequences.

The basic problem we are concerned with here is how to investigate the major statistical features of CM's. One approach studies CM's from the point of view of a random walker [2,3]. This essentially converts spatial information into time series that can be characterized, for example, by autocorrelation functions. Another approach attempts to devise tools that reflect the statistical features of CM's as a whole. Some of these features cannot be accessed with random walks alone.

In Sec. II, we generalize the autocorrelation function to the case of combinatory maps, and define density surfaces as an instance for the second approach above. Section III gives a brief classification of landscapes by autocorrelation, Sec. IV presents a study of the RNA case, Sec. V introduces Kauffman's  $n$ - $k$  model, derives landscape autocorrelation functions, and briefly reviews analytic results on gradient and adaptive walks. Section VI compares the RNA energy landscape with the  $n$ - $k$  model. Section VII concludes the paper.

## II. AUTOCORRELATION AND DENSITY SURFACES

In the case of landscapes, a random walk  $\{x_1, x_2, \dots\}$  on a configuration space generates a real-valued time series  $\{f(x_1), f(x_2), \dots\}$  whose autocorrelation function is given by

$$r(s) = \frac{\langle f(x_t)f(x_{t+s}) \rangle_t - \langle f \rangle^2}{\sigma^2}, \quad (1)$$

with variance  $\sigma^2 = \langle f^2 \rangle - \langle f \rangle^2$ . A landscape autocorre-

lation function, however, should yield information about the average changes of  $f(x_{t+s})$  as the configuration-space distance  $d$  of  $x_{t+s}$  to some reference point  $x_t$  is varied. This leads to the definition

$$\rho(d) = \frac{\langle f(x)f(y) \rangle_{d(x,y)=d} - \langle f \rangle^2}{\sigma^2}, \quad (2)$$

where  $\langle \rangle_{d(x,y)=d}$  denotes an average over all pairs of configurations  $(x, y)$  at distance  $d$  from each other.

The autocorrelation of the landscape and the autocorrelation along a walk on the landscape are related to each other via

$$r(s) = \sum_d \varphi_{sd} \rho(d), \quad (3)$$

where  $\varphi_{sd}$  denotes the probability that a walk of  $s$  steps ends at a distance  $d$  from the starting point. Equation (3) establishes a complete correspondence between the random-walk and the landscape autocorrelations, provided the  $\varphi_{sd}$  are independent of the initial conditions of the walk. For this to be the case, a sufficient regularity of the configuration space is required. Such a regularity is missing, for example, on spaces of sequences with variable length that result from insertion and deletion operations. On sequence spaces with fixed length, we obtain [4]

$$\begin{aligned} \varphi_{sd} = & \varphi_{s-1, d-1} \frac{v-d+1}{v} + \varphi_{s-1, d} \frac{d}{v} \frac{\kappa-2}{\kappa-1} \\ & + \varphi_{s-1, d+1} \frac{d+1}{v} \frac{1}{\kappa-1}. \end{aligned} \quad (4)$$

In the generic case of CM's, where the range of  $f$  could be a nonscalar object—for example, a structure—we must generalize the definition of  $\rho(d)$  with the help of the distance measure on  $y$ . We propose the following.

*Definition.*

$$\rho(d) = 1 - \frac{\langle d_y [f(x), f(y)]^2 \rangle_{d(x,y)=d}}{\langle d_y^2 \rangle}. \quad (5)$$

It is readily seen that, for landscapes  $d_y(a, b) = |a - b|$ , we recover the usual autocorrelation:

$$\begin{aligned} \langle d_y^2 \rangle - \langle d_y [f(x), f(y)]^2 \rangle_{d(x,y)=d} \\ = \langle [f(p) - f(q)]^2 \rangle - \langle [f(x) - f(y)]^2 \rangle_{d(x,y)=d} \\ = 2\langle f^2 \rangle - 2\langle f \rangle^2 - 2\langle f^2 \rangle + 2\langle f(x)f(y) \rangle_{d(x,y)=d} \\ = 2\langle f(x)f(y) \rangle_{d(x,y)=d} - 2\langle f \rangle^2. \end{aligned}$$

Here  $\langle [f(p) - f(q)]^2 \rangle = 2(\langle f^2 \rangle - \langle f \rangle^2)$  denotes the average over pairs of randomly chosen configurations  $p, q$ .

Further insight into the statistical properties of CM's is provided by density surfaces [5]. A density surface  $P(t|s)$  is the conditional probability that the images  $f(x)$  and  $f(y)$  have distance  $d_y[f(x), f(y)] = t$ , given that the configurations  $x$  and  $y$  are at a distance  $d_x[x, y] = s$  from each other. The density surface shows how, and how fast, the distribution of image differences changes as the configurations become uncorrelated. This approach provides more information than can be obtained by random walks alone. The autocorrelation  $\rho(d)$ , Eq. (5), is extract-

ed from the density surface with the following correspondences:

$$\langle d_y^2 [f(x), f(y)]^2 \rangle_{d(x,y)=s} = \sum_t t^2 P(t, |s),$$

$$\langle d_y^2 \rangle = \sum_t \sum_d t^2 P(t|s) p(s). \tag{7}$$

$p(s)$  is the frequency of configuration pairs with distance  $d_x = s$ . For sequence spaces, this amounts to

$$p(s) = \frac{1}{\kappa^n} (\kappa - 1)^s \binom{n}{s}. \tag{8}$$

From the density surface, the number of neighbors with identical image is retrieved as  $P(0|1)(\kappa - 1)n$ .

III. CLASSIFICATION BY AUTOCORRELATION

It has been suggested [6] to characterize CM's by the behavior of the autocorrelation function at small distances  $d_x$ . Landscapes with known autocorrelation function are listed in Table I.

Let the scaled distance be  $\xi = d_x / \max(d_x)$ . Many combinatorial optimization problems exhibit landscape autocorrelation functions that can be approximated for

TABLE I. Combinatorial optimization problems and their autocorrelation functions.  $\sum'$  for the  $p$ -spin model denotes the sum over all odd  $j$ , subject to the restriction  $j > \min(d, p)$ . REM is Derrida's random-energy model; STSP and ATSP denote symmetric and asymmetric traveling-salesman problems [12], GM is the graph-matching problem [13]. The corresponding metrics are transpositions (Trans.), 2-opt moves, and canonical transpositions (CTrans.). GBP is the graph bipartitioning problem [14]; its metric (Exc) is derived from exchanging a pair of objects. LAS stands for the low autocorrelated string problem [15]. The Sherrington-Kirkpatrick spin glass [16] is the special case  $p = 2$  of the  $p$ -spin model introduced in [17] as a model for a rugged landscape in evolutionary optimization. The abbreviations RN, PR, and AN refer to random neighbor, purely random, and adjacent neighbor  $N$ - $k$  model, respectively. Here the canonical metric is the Hamming metric.

Name	Metric	$\rho(d)$	$r(s)$	$l$	diam $\Gamma$	Ref.
REM	Any	$\delta_{0,d}$	$\delta_{0,s}$	0	$n$	[7]
STSP	Trans.	?	$\sim e^{-4s/n}$	$n/4$	$n - 1$	[8]
	2-opt	?	$\sim e^{-2s/n}$	$n/2$	$n/2 \dots n - 1$	[8]
	CTrans.	?	$\sim e^{-2s/n}$	$n/2$	$\frac{n(n-1)}{2}$	[6]
ATSP	Tr	?	$\sim e^{-4s/n}$	$n/4$	$n - 1$	[8]
	2-opt	?	$\sim \frac{1}{2}(\delta_{0,s} + e^{-2s/n})$			[8]
	CTrans.	?	$\sim e^{-3s/n}$	$n/3$	$\frac{n(n-1)}{2}$	[6]
GM	Trans.	?	$\sim e^{-4s/n}$	$n/4$	$n - 1$	[9]
GBP	Exc.	$1 - \frac{n-1}{n-2} \left[ 8 \frac{d}{n} - 16 \left( \frac{d}{n} \right)^2 \right]$	$\left[ 1 - \frac{8}{n} + \frac{8}{n^2} \right]^s$	$(n-3)/8$	$n/2$	[10]
LAS	Ham.	?	$\sim e^{-10s/n}$	$n/10$	$n$	[9]
RN ( $N$ - $k$ )	Ham.	$\left[ 1 - \frac{d}{n} \right] \left[ 1 - \frac{k}{n-1} \right]^d$	$\sim \left[ 1 - \frac{k+1}{n} \right]^s$	$n/(k+1)$	$n$	[11,a]
PR ( $N$ - $k$ )	Ham.	$\left[ 1 - \frac{k+1}{n} \right]^d$	$\sim \left[ 1 - \frac{k+1}{n} \right]^s$	$n/(k+1)$	$n$	[a]
AN ( $N$ - $k$ )	Ham.	Eq. (24c)	$\sim \left[ 1 - \frac{k+1}{n} \right]^s$	$n/(k+1)$	$n$	[11,a]
$p$ -spin	Ham.	$1 - \frac{2}{\binom{n}{p}} \sum' \binom{d}{j} \binom{n-d}{p-j}$	$\sim e^{-2ps/n}$	$n/(2p)$	$n$	[6]
SK	Ham.	$1 - \frac{n}{n-1} \left[ 4 \frac{d}{n} - 4 \left( \frac{d}{n} \right)^2 \right]$	$\left[ 1 - \frac{4}{n} \right]^s$	$n/4$	$n$	[6]

<sup>a</sup>This paper.

small  $\xi$  by the first terms in the expansion of the exponential,

$$\rho(\xi) = 1 - \frac{1}{\alpha} \xi + \dots \quad (9)$$

The examples of Table I all induce random walks with exponential autocorrelation functions  $r(s)$ . This need not be generally the case, as exemplified by superpositions of independent landscapes with different characteristic lengths that obviously result in landscapes with no unique  $\alpha$ . Even if a walk on a landscape does not exhibit exactly a decaying exponential as autocorrelation function, a rough correlation indicator is given by the nearest-neighbor correlation:

$$\lambda = n\alpha \approx \frac{\langle d_y^2 \rangle}{\langle d_y [f(x), f(y)]^2 \rangle_{d(x,y)=1}} \quad (10)$$

Another class of landscapes is characterized by discontinuous autocorrelation functions, as, for example, the random energy model [7], or the asymmetric traveling salesman problem (TSP) [8]; see Table I.

Smooth landscapes have  $\rho(\xi) = 1 - \xi^2 + \dots$ , and non-trivial fractal landscapes are characterized by  $\rho(\xi) = 1 - |\xi|^\alpha + \dots$ , with  $\alpha \neq 0, 1, 2$ .

#### IV. RNA LANDSCAPES

##### A. RNA folding

An RNA sequence of length  $n$  is represented as a string  $I = [s_1 s_2 \dots s_n]$ , with the  $s_i$  taken from an alphabet  $\mathcal{A}$ . Here we consider the natural four-letter alphabet  $\mathcal{A} = \{A, U, G, C\}$ , the binary alphabets  $\mathcal{A} = \{G, C\}$  and  $\mathcal{A} = \{A, U\}$ , as well as artificial four- and six-letter alphabets  $\mathcal{A} = \{A, B, C, D\}$  and  $\mathcal{A} = \{A, B, C, D, E, F\}$ , respectively, where complementarity is assumed between A and B, C and D, E and F (with energy parameters as in the GC case).

RNA structure can be broken down conceptually into a secondary structure and a tertiary structure. The secondary structure is a pattern of complementary base pairings (Fig. 1). The tertiary structure is the three-dimensional (3D) configuration of the molecule. As opposed to the protein case, the secondary structure of RNA sequences is well defined; it provides the major set of distance constraints that guide the formation of tertiary structure, and covers the dominant energy contribution to the 3D structure. In this paper, we will be concerned only with secondary structures.

The definition of secondary structure used in most computational approaches assumes planarity. Planarity means that unpaired nucleotides inside a loop cannot pair with unpaired nucleotides outside a loop. It is known, however, that unpaired bases from different loop regions may pair with each other, forming so-called pseudoknots. While the computational problem for strictly planar secondary structures has been essentially solved in the early 1980s [18–21], the problem involving long-range pseudoknots is still unsolved.

A secondary structure  $S(I) \in \mathcal{S}$  is formally defined as the set of all base pairs  $(s_i, s_j)$  with  $(i < j)$  fulfilling two

requirements [19]: (1) each base is involved in at most one base pair, and (2) there are no knots or pseudoknots, i.e., if  $(s_i, s_j)$  and  $(s_k, s_l)$  are base pairs, then  $i < k < l < j$  or  $k < i < j < l$ .

The basic elements of secondary structures are shown in Fig. 1. For the free energies of these building blocks, experimental data are available. The elements are as-

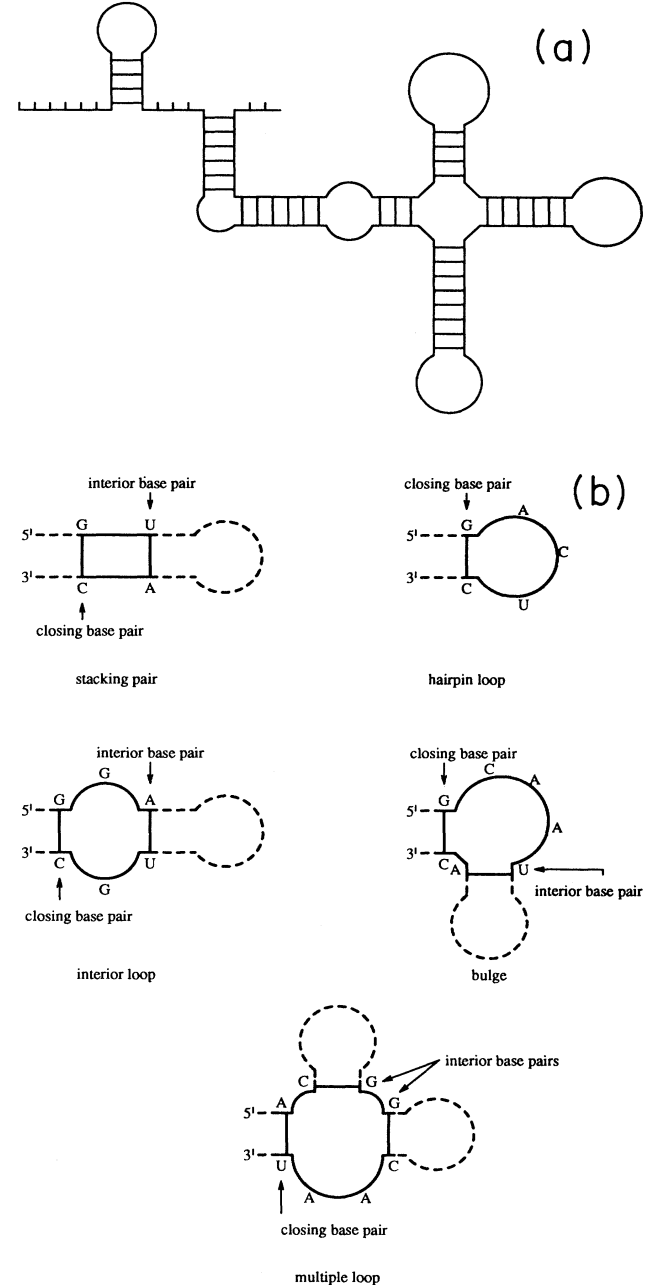


FIG. 1. Top: example of a secondary structure. Below: basic structure elements. Every secondary structure can be decomposed into such basic elements. The free energy of a secondary structure is the sum over all contributions of these elements. Their energy contributions have been experimentally determined as a function of the nucleotide sequence [22–24].

sumed to contribute additively to the overall free energy of the complete secondary structure. The data set used here has been taken from the literature [22–24].

Under thermodynamically reasonable assumptions, the folding problem for planar RNA secondary structure consists in finding a minimum free-energy structure, and can be attacked by the technique of dynamic programming. Our implementation of the folding algorithm follows the reasoning given by Zuker and Stiegler [20]. The combinatory map from the space of sequences  $\mathcal{A}^n$ , into the space of minimum free-energy structure  $\mathcal{S}$ ,

$$S: \mathcal{A}^n \rightarrow \mathcal{S}, \quad I \rightarrow S(I), \quad (11)$$

is the quintuple  $(\mathcal{A}^n, d_H, \mathcal{S}, d_{\mathcal{S}}, S)$ , with  $d_H$  denoting the Hamming distance, and  $d_{\mathcal{S}}$  being a suitably defined distance between secondary structures (see below). If we view the folding procedure as a map from a sequence to the free energy of the corresponding structure,

$$\Delta G: \mathcal{A}^n \rightarrow \mathbb{R}, \quad I \rightarrow \Delta G(I), \quad (12)$$

we obtain the combinatory landscape  $(\mathcal{A}^n, d_H, \mathbb{R}, |\cdot \cdot \cdot|, \Delta G)$ .

### B. Distribution of energies

Here we focus on global features relating to the distribution of free-energy values over the RNA landscape. Table II shows the mean free-energy values  $\langle \Delta G \rangle$  and

variances  $\sigma^2$  for a variety of alphabets. We have also computed the skewness (third moment scaled by  $\sigma^3$ ) and the kurtosis (fourth moment scaled by  $\sigma^4$ ) of the distribution (not shown).

The main observations are that  $\langle \Delta G \rangle$  scales linearly with chain length  $n$  for all alphabets, and so does the variance  $\sigma^2$  for  $n > 50$  (see Table IV below). The distribution sharpens considerably for longer chains. A Gaussian distribution is characterized by vanishing skewness and a value of 3 for the kurtosis. This seems to be mostly the case for the biophysical GCAU alphabet and long chains. The other alphabets show deviations from the Gaussian case in skewness and/or kurtosis. The available data do not allow a definite statement about the limiting behavior for very large  $n$  to be made. We did not pursue this issue in depth because of doubts about the validity of a thermodynamic (rather than kinetic) folding algorithm in that limit.

### C. Free-energy autocorrelation

Figure 4(a) shows an example of the landscape autocorrelation function  $\rho(d)$ . Complementary sequences have similar structures and energies. For each reference sequence on binary alphabets, there is only one complement, and, therefore, this approximate symmetry shows up as a U-shaped correlation. Table III lists the characteristic lengths of the landscape, Eq. (10), for various al-

TABLE II. Averages of the free energies  $\Delta G$  and their variances  $\sigma^2$  for (1) GC, (2) AU, (3) GCAU, (4) GCXK, and (5) ABCDEF. For (3), the biophysical parameter set is used; for (4) and (5), the GC parameter set is used for all base pairs. Energy unit: 0.1 kcal mol<sup>-1</sup>. Energy data set taken from [23].

$n$	$-\Delta G$					$\sigma$				
	(1)	(2)	(3)	(4)	(5)	(1)	(2)	(3)	(4)	(5)
20	74.40	2.22	10.29	19.06	6.59	28.80	5.12	14.63	19.95	12.35
25	108.85	5.68	17.56	31.50	16.05	31.86	8.78	18.56	23.88	16.05
28				39.38					25.79	
30	144.23	9.59	25.58	44.46	16.91	34.55	11.24	21.81	26.96	18.89
30 <sup>a</sup>			25.68					22.00		
35	180.27	14.09	33.61	58.24	22.71	36.96	13.18	24.35	29.22	21.32
40	217.00	18.83	42.33	72.60	29.30	39.17	14.64	26.88	31.54	23.59
45	253.73	23.53	50.58	86.98	35.67	40.87	15.81	28.83	33.25	25.38
50	290.61	28.54	60.12	101.89	42.36	42.66	16.89	31.03	35.32	27.24
50 <sup>a</sup>			60.16					31.13		
55	327.59	33.54	68.84	116.47	48.89	44.23	17.81	32.82	36.89	28.70
60	364.50	38.61	77.48	131.28	55.81	45.51	18.79	34.63	38.19	30.10
70	440.73	48.46	95.64	161.63	69.59	47.95	20.25	37.64	40.99	32.66
80	516.83	59.32	113.77	192.19	83.57	50.76	21.78	40.46	43.37	34.87
100	669.35	80.12	150.78	254.56	111.46	54.21	24.17	45.96	48.04	38.90
120	824.41					58.42				
150	1054.77	133.45	244.57			62.12	28.79	57.06		
200	1445.62	188.21	339.16			68.63	32.58	66.28		
250	1838.36	243.62	434.93			74.08	36.04	73.79		
300	2233.63	299.62	531.34			79.80	39.70	81.14		
350	2629.50	355.88	628.59			84.40	42.72	87.96		
400	3026.54	412.50	726.02			88.44	45.50	94.41		
450	3423.91	469.30	823.52			92.98	47.70	99.92		
500	3822.34	526.24	921.46			95.74	49.69	105.2		

<sup>a</sup>Specially stable tetraloops taken into account.



### D. Structure autocorrelation

As detailed in Sec. II, computing the autocorrelation of the structures themselves requires a distance measure on the space of structures. The definition is essentially based on converting one-to-one secondary structure graphs into rooted ordered trees [5,26–28]. The distance between two trees is obtained as a generalization of sequence alignment procedures, and involves minimizing the cost of transforming one tree into the other with elementary editing operations, such as deletion, insertion, and relabeling of nodes. For details, see [5].

Figure 4 shows the structure autocorrelation function  $\rho(d)$ . The U-shaped form for the binary alphabet arises

for the same reason given in Sec. IV C.

The structure density surfaces of sequences with  $n = 100$  for the GC and GCAU alphabets are shown in Fig. 5. First, we note that nearest neighbors in configuration space can exhibit already substantially different structures with significant probability. Furthermore, it is extremely unlikely that two randomly chosen sequences fold into identical structures. This is in sharp contrast to the energy landscape. Stated differently: there are many more structures than energies. Nonetheless, the distribution of structure distances approaches already, at fairly short Hamming distances, the distribution expected for random sequences. In contrast to the free energy, there is no size-independent upper bound to the

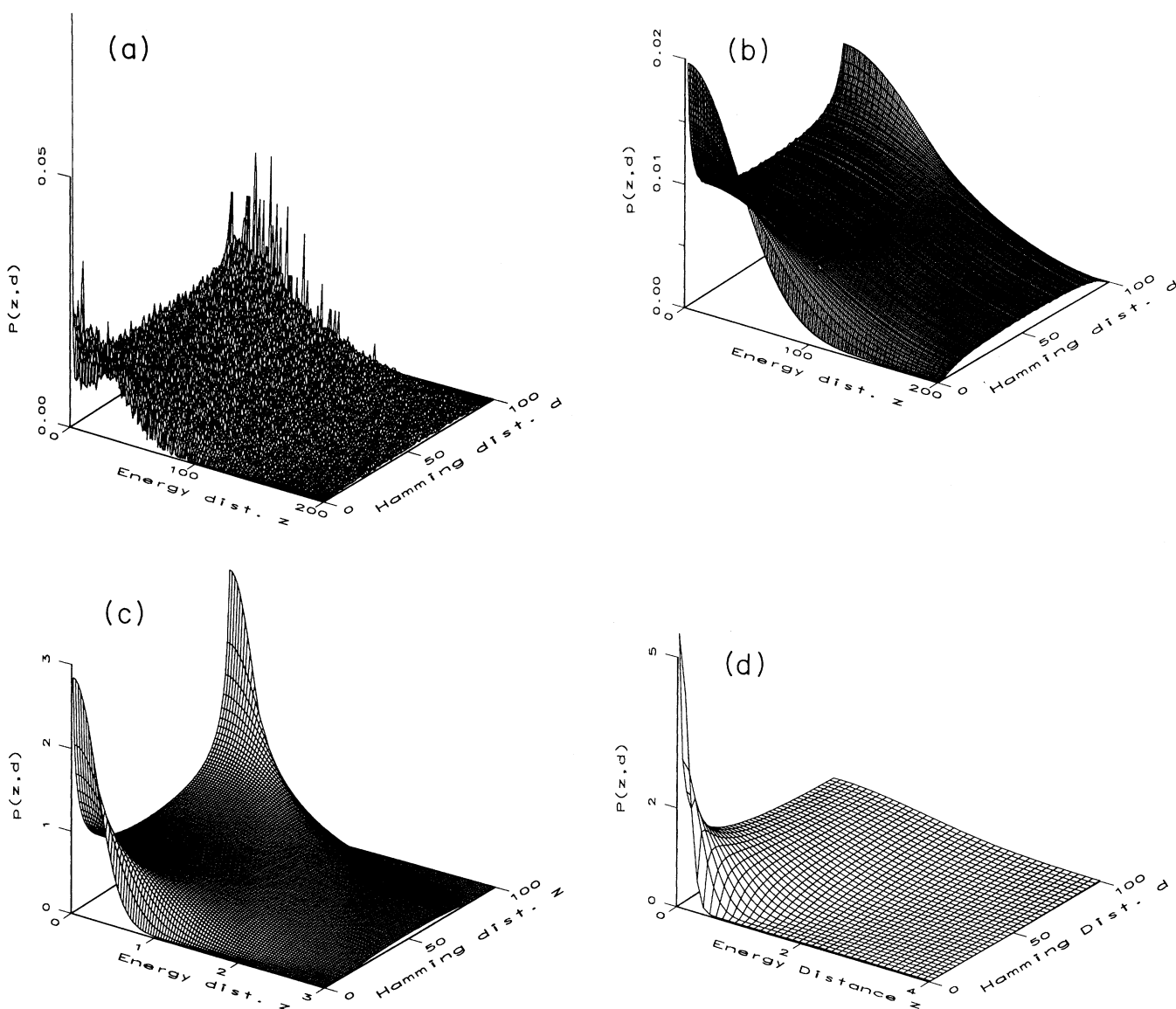


FIG. 2. (a) Sampled density surface  $P(z|d)$  for energy differences of minimum-free-energy structures on GC sequences of length  $n = 100$ . (b) Gaussian approximation with the statistically determined autocorrelation function  $\rho(d)$ . (c) Exact density surface for the Sherrington-Kirkpatrick spin glass. (d) Exact density surface for a RN  $n-k$  model with  $k = 7$ .

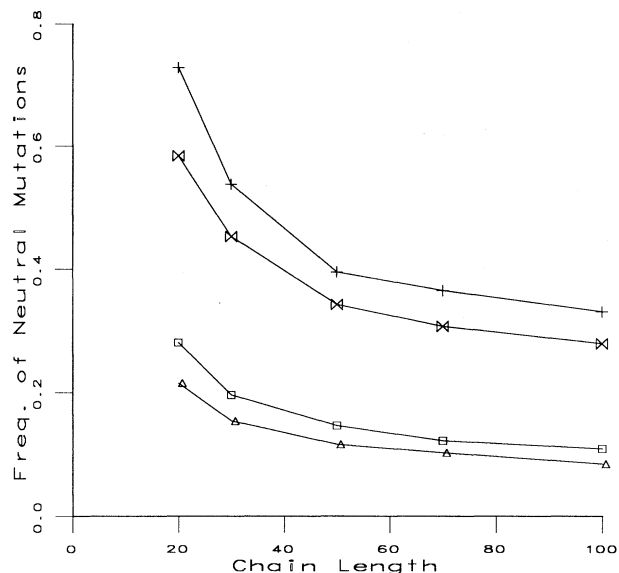


FIG. 3. Frequency of neutral neighbors  $P(0|1)$  on the free-energy landscape as a function of chain length. The upper two curves refer to GCAU sequences, the lower two curves refer to GC sequences. Data points designated with crosses and squares were obtained with energy parameters taken from [22], while triangles and bowties refer to energy parameters taken from [23].

distance between secondary structures of two neighboring sequences. This suggests that, in a small ball around a random sequence, one may find the vast majority of all minimum-free-energy structures  $\mathcal{S}$  [5]. The phenomenon is known as “shape space covering” [29]. We will report elsewhere on the verification of this conjecture with additional methods. Although there are many fewer different energy values than minimum energy structures,

TABLE V. Characteristic lengths of RNA secondary structure landscapes for various alphabets and chain lengths. Energy parameters are taken from [23].

$n$	GC	AU	GCAU	GCXK	ABCDEF
20	1.51	1.22	2.84	2.43	2.50
25	1.69	1.61	3.68	2.88	3.28
30	1.87	1.85	3.99	3.25	3.92
35	2.02	2.12	4.49	3.57	4.51
40	2.15	2.24	4.85	3.74	
45	2.33	2.33	4.97	4.10	5.49
50	2.35	2.50	5.46	4.33	5.84
55	2.62	2.64			6.20
60	2.68	2.83	6.01	4.84	6.60
70	2.87	2.91	6.25	5.24	7.14
80	3.13	3.32		5.67	7.93
90		3.60			
100	3.36	3.77	7.63	6.45	8.94
120	3.72				

the energies are not scattered randomly across the sequence space.

The characteristic lengths of the structure autocorrelation (Table V) are shorter by a factor between 2 and 3 compared to their free-energy counterparts. They roughly correspond to the Hamming distance at which the dominant peak resulting from identical or very similar structures has dissolved. Randomization of structure for GCAU sequences of length  $n=100$  occurs already at a distance of about  $d_H=16$  from a random reference point (Fig. 5). From Table V, we see that this corresponds to a ball of radius  $d_H \approx 2\lambda$ . The characteristic length data are roughly consistent with linear scaling in  $n$ .

The ordering of characteristic lengths of the structure autocorrelation for long sequences is summarized by

$$[\text{AU}] \approx [\text{GC}] \ll [\text{GCXK}] < [\text{GCAU}] < [\text{ABCDEF}] . \quad (15)$$

TABLE IV. Scaling with sequence length of average free energy, variance, and correlation length of free energies for different alphabets. Energy parameters from [23].

	GC	AU	GCAU	GCXK	ABCDEF
			$-\langle \Delta G \rangle$		
Slope	7.8120	1.0948	1.9006	2.9495	1.3044
Interp.	-101.21	-26.14	-32.25	-43.97	-21.40
Corr.	0.99996	0.9998	0.99992	0.9996	0.998
			$\sigma^2$		
Slope	17.33	5.06	22.61	23.56	17.07
Interp.	959.8	22.2	-180.4	16.1	-143.3
Corr.	0.996	0.997	0.99996	0.997	0.998
			$\lambda$		
Slope	0.0857	0.0600	0.2627	0.1078	0.1182
Interp.	0.43	0.06	-1.93	0.50	0.47
Corr.	0.992	0.997	0.997	0.993	0.997



Binary alphabets generally form more structures because the probability for two randomly chosen positions along the sequence to pair is highest. Changing one position, therefore, is more likely to alter the minimum energy structure. In contrast, the characteristic lengths of the free-energy autocorrelation do not follow the same ordering:

$$[\text{AU}] < [\text{GC}] < [\text{GCXK}] < [\text{ABCDEF}] \ll [\text{GCAU}]. \quad (16)$$

### E. Biased walks and local optima

Additional information about the local structure of landscapes is provided by biased walks. These are random walks aimed at reaching local extrema of the landscape. The data provided by this technique will be

compared with a statistical model in the following sections.

We define a local minimum  $\hat{y}$  in configuration space by

$$f(\hat{y}) \leq f(x) \text{ for all neighbors } x \text{ of } \hat{y}. \quad (17)$$

Local maxima are defined analogously. The term local optimum means either minimum or maximum, according to context.

We consider here two types of biased walks, adaptive walks and gradient walks. In an adaptive walk, a starting point  $x_0$  is chosen at random. The walk then proceeds to a randomly chosen neighbor  $x'$  such that  $f(x') < f(x_0)$ , and it stops if no neighbor satisfies this condition, at which point the walk has reached a local minimum. A gradient walk, by contrast, proceeds to the neighbor  $x'$  for which  $f(x')$  is minimal.

Table VI compiles the average lengths of adaptive and

TABLE VI. Adaptive and gradient walks on RNA free-energy landscapes.  $\bar{F}$  denotes the free energy of random sequences. Average free energies of local optima are denoted by  $\bar{F}_{\text{opt}}$ ,  $\bar{F}_{\text{opt}}^*$ , and  $\bar{F}_{\text{opt}}^{**}$ , depending on whether they are generated randomly, as end points of adaptive walks, or as end points of gradient walks, respectively. The standard deviations of their distributions are denoted by  $\sigma_F$ ,  $\sigma_{\text{opt}}$ ,  $\sigma_{\text{opt}}^*$ , and  $\sigma_{\text{opt}}^{**}$ . The lengths of adaptive and gradient walks are denoted by  $L_{\text{adap}}$  and  $L_{\text{grad}}$ , respectively. The rightmost column compiles the probability  $p_{\text{LO}}$  for finding a local optimum at random.

$n$	$\bar{F}$	$\sigma_F$	$\bar{F}_{\text{opt}}$	$\sigma_{\text{opt}}$	$\bar{F}_{\text{opt}}^*$	$\sigma_{\text{opt}}^*$	$\bar{F}_{\text{opt}}^{**}$	$\sigma_{\text{opt}}^{**}$	$\lambda$	$L_{\text{adap}}$	$L_{\text{grad}}$	$p_{\text{LO}}$
GC												
Energy parameters from [22]												
20	118.9	43.0	179.7	32.0	204.5	29.1	211.5	31.2	2.10	2.42	1.72	0.1138
25	178.9	47.3	257.6	44.1	300.4	37.8	305.2	42.0	2.70	3.30	2.29	0.0530
30	239.7	50.8	329.3	50.6	398.5	48.7	402.9	53.7	3.09	4.44	3.12	0.0220
35	303.1	54.0	390.3	43.6	493.1	60.5	496.5	67.4	3.57	5.49	3.88	0.0108
40	368.6	57.6	465.5	48.8	584.7	74.5	585.3	79.1	3.98	6.42	4.48	0.0084
50	500.1	62.3	633.4	48.9	759.8	90.7	755.5	92.4	4.67	8.06	5.50	0.0025
60	632.2	66.7	768.0		934.4	97.4	930.9	94.9	5.65	9.61	6.42	0.0008
Energy parameters from [23]												
20	74.4	28.8	127.9	24.9	148.3	21.6	152.0	22.2	1.96	3.55	2.51	0.0230
30	144.2	34.6	213.1	28.1	262.6	36.6	265.1	37.3	2.93	5.32	3.70	0.0050
40	217.0	39.2	315.0	29.5	369.2	44.4	372.4	44.7	4.00	6.78	4.69	0.0011
50	290.6	42.7	415.0		478.0	48.9	481.9	49.3	4.51	8.27	5.68	0.0002
60	365.5	45.5			587.6	56.6	589.4	55.1	6.17	9.62	6.64	
GCAU												
Energy parameters from [22]												
30	34.2	31.0	94.6	45.4	208.5	72.3	228.5	78.4	6.73	8.01	5.12	0.0009
40	59.8	37.9			307.3	93.5	339.8	103.7	9.85	11.62	7.57	
50	90.6	43.8			404.1	105.1	441.8	112.0	13.76	14.77	9.31	
60	128.1	48.1			504.6	118.0			18.38	17.99		
70	162.7	50.6			604.2	125.3	644.1	132.9	20.80	21.22	13.32	
80	201.9	55.3			704.2	139.2			23.81	24.17		
Energy parameters from [23]												
30	25.7	22.0			172.9	43.9	181.1	47.4	6.08	11.84	7.17	
40	42.3	26.9			253.0	50.4	272.0	54.8	8.14	16.26	9.95	
50	60.2	31.1			325.8	59.6	352.3	65.1	11.74	20.73	12.53	
60	77.5	34.6			401.8	68.3	428.0	71.9	13.20	25.30	14.92	
70	95.6	37.6			478.4	73.1	508.0	80.1	16.20	29.13	17.44	

gradient walks on the free-energy landscape, as well as the average free energy (and the corresponding variances) of start and end points. These quantities were calculated for an energy parameter set available in the late 1970s [22] and for a recently updated parameter set [23,24].

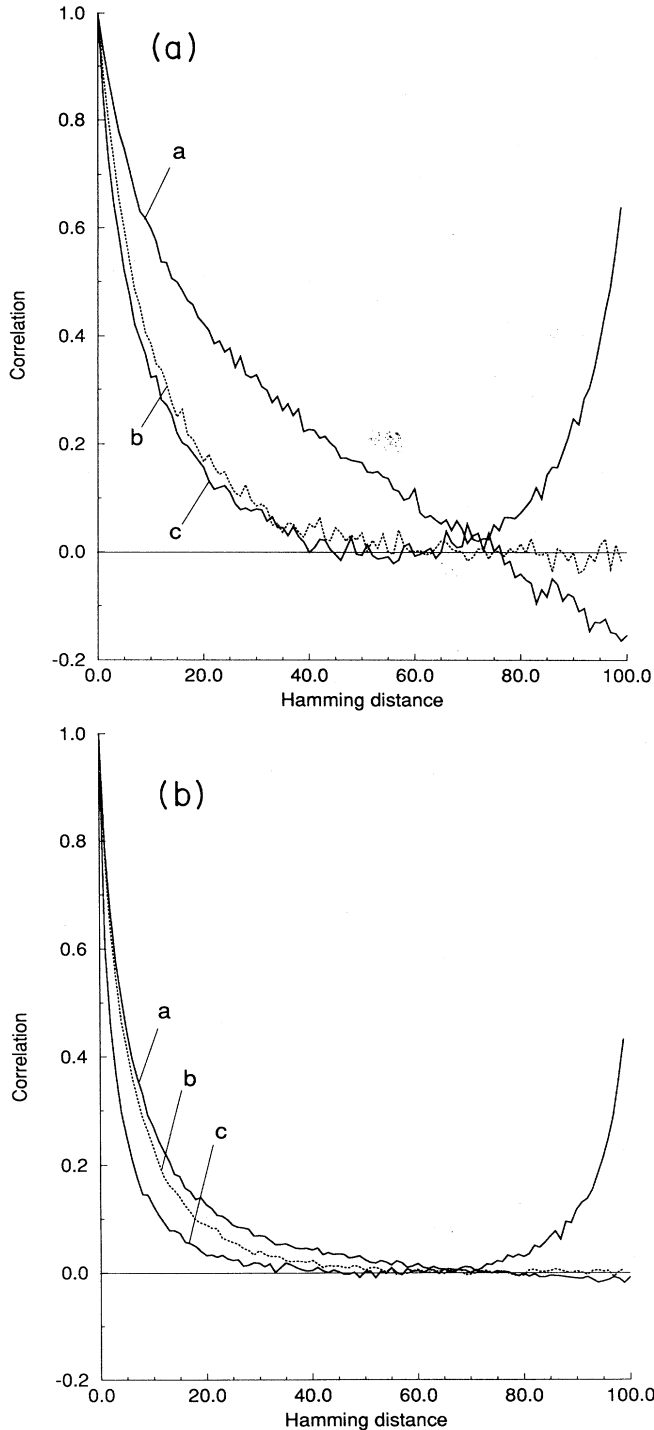


FIG. 4. Autocorrelation functions for sequences of length  $n = 100$  for free-energy landscapes (a), and (b) secondary structure landscape, with alphabets GCAU (a), GCXK (b), and GC (c).

While the absolute values differ, the scaling behavior is not affected. The average free energy of local optima reached by adaptive and gradient walks is roughly the same, and scales linearly with  $n$ . The rightmost column in Table VI summarizes the frequency  $p_{LO}$  with which the initially chosen sequence was already a local optimum. We refer to such a sequence as a “random optimum.” The frequency  $p_{LO}$  decreases exponentially with chain length  $n$  (see Sec. IV E). This frequency was, however, too low to yield meaningful quantities for GCAU sequences in both parameter sets. We note that the energy of local optima reached at the end of the biased walks is substantially lower than the average free energy of random optima.

The distribution of free-energy increments per step along a random walk is shown in Fig. 6. The local deviations from a Gaussian distribution are again most pronounced at 0, resulting from the substantial fraction of neutral neighbors. In the case of binary sequences [Fig. 6(a)], particular free-energy changes are favored. As a biased walk proceeds, the distribution of energy increments is obviously altered. After some five steps, the dis-

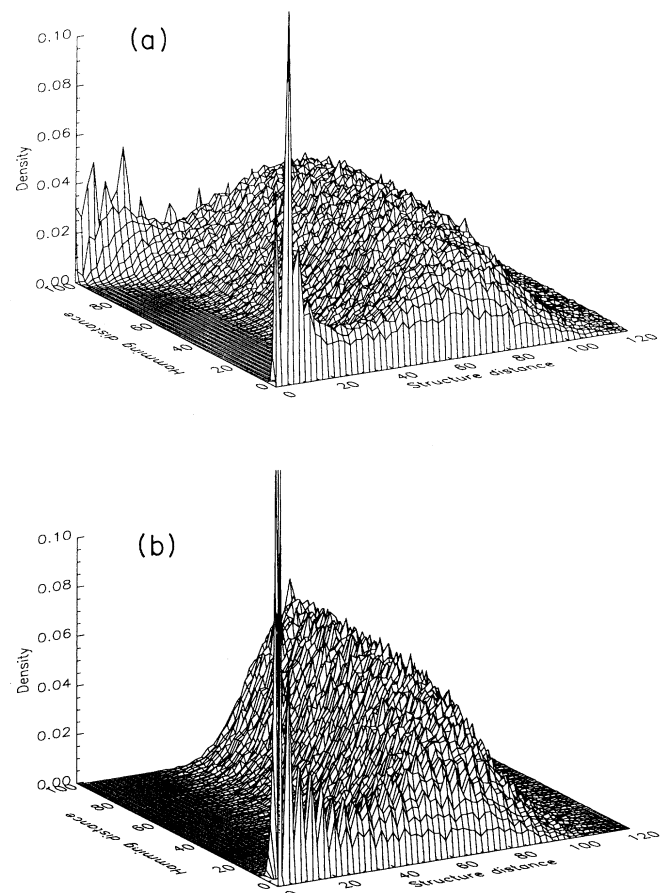


FIG. 5. Density surface for secondary structure differences as determined by tree editing. Parameters are as in Fig. 2.  $n = 100$ ; upper part: GC, lower part: GCAU.

tribution is dominated by a few quanta (Fig. 7) that identify the most likely structural changes.

In Fig. 8, we plot the average change in free energy for landscapes of different length along biased walks. The number of walk steps is scaled by the corresponding characteristic length  $\lambda$ . The figure demonstrates that for a fixed alphabet the properly scaled statistical features of walks are independent of the system size.

### V. KAUFFMAN'S $n$ - $k$ MODEL

Because of the important role of landscapes in problems of optimization and evolutionary adaptation [30],

there has been considerable interest in devising simple statistical models, the hope being that these models reflect the proper qualitative features of their natural counterparts. In Sec. IV we had a close look at an expensive but realistic model of a landscape induced by RNA folding. In this section, we present a frequently used statistical model due to S. Kauffman, the so-called  $n$ - $k$  model [31,32].

Let us assume that the energy  $F$  of a string of  $n$  bits is the average of contributions from each of the individual bits. We choose the contribution from the  $i$ th bit,  $f_i$ , as a random function of the state of that bit and a context given by  $k < n$  other bits. Each of the  $2^{k+1}$  possible

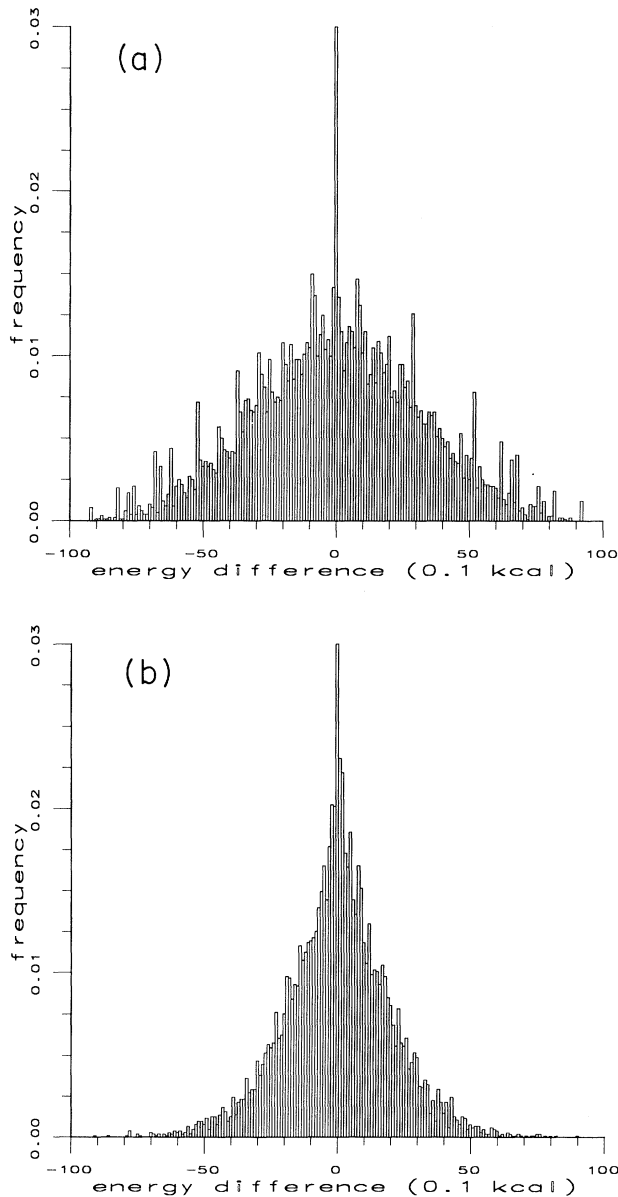


FIG. 6. Distribution of per-step energy differences along a random walk. (a) GC sequences of length  $n=100$ . (b) AUGC sequences of length  $n=100$ .

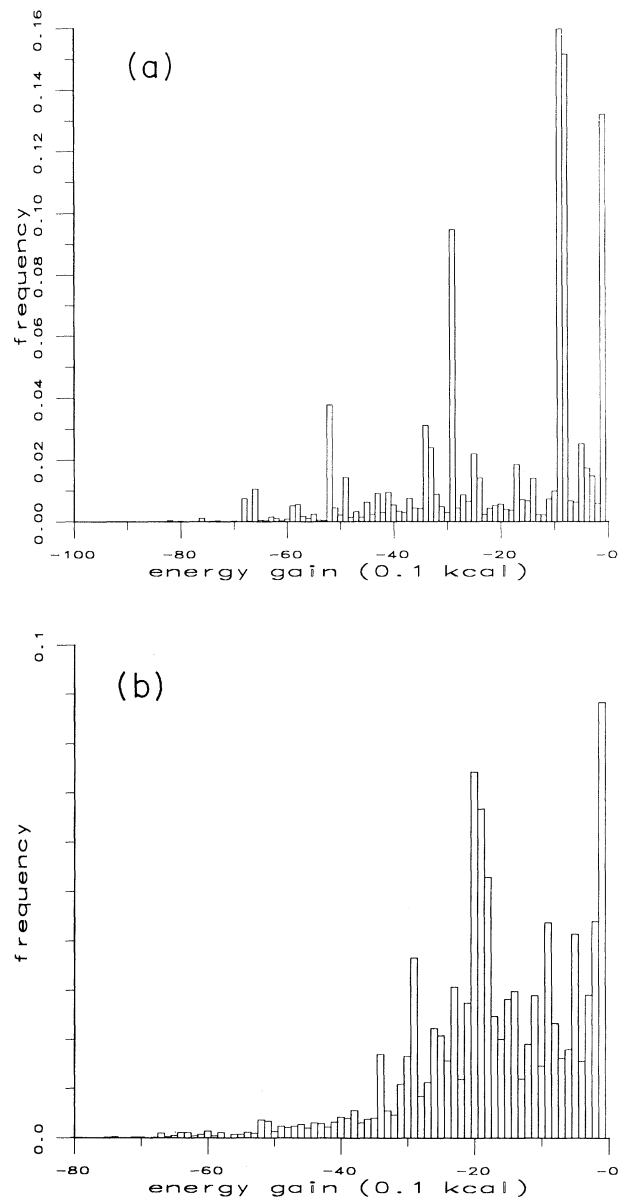


FIG. 7. Distribution of energy increments after more than four steps along a gradient walk. (a) GC landscape and (b) AUGC landscape.

values for  $f_i$ , one for each of the possible states of the  $k + 1$  bits upon which  $f_i$  depends, is assigned by selecting an independent random variable from some specified probability distribution,  $P(x) = \Pr(f_i \leq x)$ . This set of assignments constitutes the “energy table” for the  $i$ th bit. There is a different, independently generated table for each of the  $n$  bits, which, once chosen, is never reassigned.

It remains to specify the  $k$  sites that influence a given

position. The simplest choice is to imagine that the neighborhoods consist of the  $k/2$  bits that are the nearest neighbors on each side of the bit in question and to assume that the bits are arranged in a circle (periodic boundary conditions). The other extreme is to pick the  $k$  sites at random. In a variant of this model, it is not required that  $f_i$  depends on  $i$ , but all  $k + 1$  relevant sites are chosen randomly. These three versions of the  $n$ - $k$  model will be referred to as AN (adjacent neighborhood), RN (random neighborhood), and PR (purely random) model, respectively. These two extremes correspond to two important types of spin glasses: the adjacent neighborhoods correspond to a one-dimensional, short-range spin glass; the random neighborhoods correspond to a dilute, long-range spin glass.

If the underlying probability distribution of the site energies,  $f_i$ , has a finite variance, the distribution of the fitness of the entire string,

$$F = \frac{1}{n} \sum_i f_i, \quad (18)$$

will tend to a Gaussian with mean  $\mu$  and variance  $\sigma^2/n$ , where  $\mu$  and  $\sigma^2$  are, respectively, the mean and variance of the  $f_i$ 's as  $n \rightarrow \infty$ .

#### A. Correlation

Weinberger [3] shows that the autocorrelation function of a random walk on the  $n$ - $k$  model landscape is a single decaying exponential to within an error of  $O(1/n)$ . All isotropic landscapes of this kind can be generated by an appropriate choice of the mean and variance of the site-energy distribution  $P(x)$ , and an appropriate choice of  $k$ .

The ruggedness of the landscape varies dramatically as  $k$  varies from 0 through  $n - 1$ . For  $k = 0$ , each site is independent of all other sites. The autocorrelation function in this special case reads

$$\rho(d) = 1 - d/n. \quad (19)$$

After some algebra, one shows that a random walk on this landscape generates an AR (1) process with autocorrelation function

$$r(s) = \sum_{d=0}^s \varphi(s,d) \rho(d) = \left[ 1 - \frac{1}{n} \frac{\kappa}{\kappa - 1} \right]^s. \quad (20)$$

In contrast, the  $k = n - 1$  landscape is the random-energy model: the energy contribution of each site then depends on all of the other sites because the context for each of the  $n - 1$  other bits is changed when even a single bit is flipped. In this case, therefore, the energy of each  $n$ -bit string is assigned an energy that is statistically independent of its neighbors.

The autocorrelation functions for the various types of  $n$ - $k$  models are obtained from

$$\begin{aligned} [F(x) - F(y)]^2 &= \sum_{i=1}^N [f_i(x) - f_i(y)]^2 \\ &+ 2 \sum_{i < j} [f_i(x) - f_i(y)] \\ &\quad \times [f_j(x) - f_j(y)]. \end{aligned} \quad (21)$$

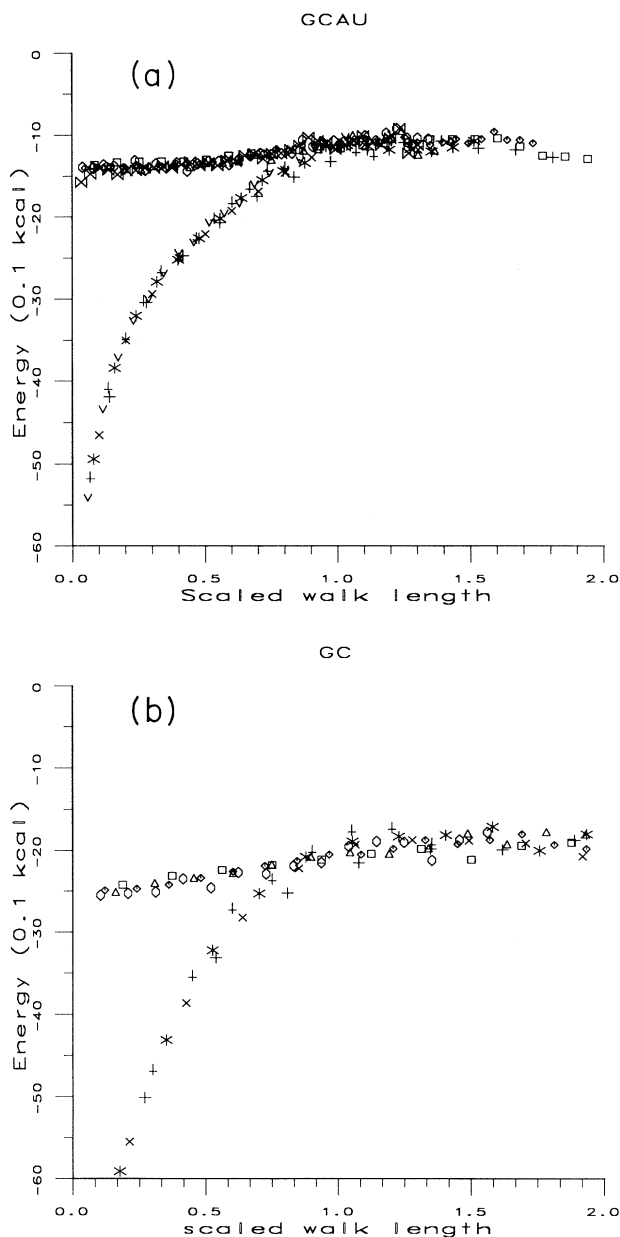


FIG. 8. Average energy increment per step for biased walks on the (a) GC landscape and (b) AUGC landscape.  $\square$ ,  $\Delta$ ,  $\diamond$ ,  $\circ$ ,  $\boxtimes$  belong to adaptive walks and chain lengths 30, 40, 50, 60, respectively. The corresponding gradient walks are denoted by  $+$ ,  $\times$ ,  $*$ ,  $\dagger$ , and  $\vee$ . Energy parameter set taken from [23].

The average over the second sum vanishes because of the independence of the  $f_i$ 's, and the contributions of the terms  $[f_i(x) - f_i(y)]^2$  are zero if  $x$  and  $y$  have the same digit at all positions that are in the neighborhood of site  $i$ . Otherwise the contributions are random with mean value  $\sigma^2$ . Thus we obtain

$$\langle [F(x) - F(y)]^2 \rangle = \frac{1}{n} 2\sigma^2 n h(d), \quad (22)$$

where  $h(d)$  is the probability for a site to appear at least once in the list of the  $d$  sites that are changed by moving from  $x$  to  $y$ . For the  $n-k$  model, it is evident that complementary strings are uncorrelated, because they do not have any contribution to site fitnesses in common; therefore the variance of  $F$  is  $\sigma^2$  and

$$\rho(d) = 1 - h(d). \quad (23)$$

For the three variants of the  $n-k$  model, the functions  $h(d)$  are derived in the Appendix:

$$\begin{aligned} h_{\text{RN}}(d) &= 1 - \left[ 1 - \frac{d}{n} \right] \left[ 1 - \frac{k}{n-1} \right]^d, \\ h_{\text{PR}}(d) &= 1 - \left[ 1 - \frac{k+1}{n} \right]^d \\ h_{\text{AN}}(d) &= \frac{k+1}{n} d - \frac{1}{\binom{d}{n}} \\ &\quad \times \sum_{j=1}^{\min(k, n+1-d)} (k-j+1) \binom{n-j-1}{d-2} \end{aligned} \quad (24)$$

for  $d \geq 2$ .

Evidently,  $h(d)$  is independent of the number of letters in the alphabet. For all three types of  $n-k$  models, the correlation length can be estimated from the nearest-neighbor correlations, Eq. (10),

$$\rho(1) = \left[ 1 - \frac{1}{n} \right] \left[ 1 - \frac{k}{n-1} \right] = 1 - \frac{k+1}{n}, \quad (25)$$

since all  $h(1)$  are identical. This yields the estimate

$$\lambda = -1 / \ln \left[ 1 - \frac{k+1}{n} \right] = \frac{n}{k+1} - \frac{1}{2} + O \left[ \frac{k+1}{n} \right]. \quad (26)$$

### B. Biased walks and local optima

In the case  $k=0$ , there is, with probability 1, a unique optimal digit for each site; hence, a single specific sequence comprised of the optimal digit values in each position is almost surely the single, global optimum in the energy landscape. Any other string is suboptimal, and lies on a connected walk via nearest-neighbor better variants to the global optimum. The length of the walk is just the Hamming distance from the initial string to the global optimum. For a randomly chosen initial string,  $(\kappa-1)/\kappa$  of the digits will be in an energetically less favorable state ( $\kappa$  is the size of the alphabet), hence the expected walk length is just  $[(\kappa-1)/\kappa]n$ .

For the random-energy model, it has been shown [1,33–35] that the landscapes have many local optima, that the walks to optima are short [ $O(\ln n)$ ], and that only a small fraction of local optima are accessible from any initial string. Weinberger [36] derived expressions for the expected number of local optima and the probability distribution of their energies. In the Gaussian case, the expected number of local optima  $\mathcal{N}_{\text{LO}}$  is given by

$$\ln(\mathcal{N}_{\text{LO}}) \approx c_0 + n \left[ \ln \kappa - \frac{\ln[(\kappa-1)(k+1)]}{k+1} \right], \quad (27)$$

where  $c_0$  is some small constant. The distribution of energies of local optima is Gaussian with mean  $\langle F_{\text{LO}} \rangle = \mu - \eta\sigma$ , where

$$\eta = \left[ \frac{2 \ln[(\kappa-1)(k+1)]}{k+1} \right]^{1/2}, \quad (28)$$

and variance

$$\sigma_{\text{LO}}^2 = \sigma^2 \frac{1}{n} \frac{1}{1 + 4(\kappa-1)(k+2)[\ln(k+1)]/(k+1)}. \quad (29)$$

Estimates for the average lengths of biased walks have been derived recently [36]. The average length of gradient walks can be estimated from the average distance of two local optima as

$$L_{\text{grad}} = n \left[ \frac{\kappa-1}{\kappa} \right] \frac{\ln[(\kappa-1)(k+1)]}{(k+1) \ln \kappa}. \quad (30)$$

As an estimate for the length of adaptive walks, one obtains

$$L_{\text{adap}} = n \frac{\ln[(\kappa-1)(k+1)]}{(k+1)}. \quad (31)$$

The lengths of adaptive walks and gradient walks differ, therefore, by a constant factor

$$\frac{L_{\text{adap}}}{L_{\text{grad}}} = \frac{\kappa \ln \kappa}{\kappa - 1}. \quad (32)$$

The chance to find a local optimum at random is given by

$$\ln p_{\text{LO}} = \ln \mathcal{N}_{\text{LO}} - n \ln \kappa = c_0 - n \frac{\ln[(\kappa-1)(k+1)]}{k+1}. \quad (33)$$

## VI. COMPARISON OF THE RNA LANDSCAPE WITH THE $n-k$ MODEL

In order to compare the  $n-k$  model and the RNA landscape, we need a hypothetical  $k$  value. From Eq. (26), we obtain

$$k = n(1 - e^{-1/\lambda}) - 1. \quad (34)$$

Since the characteristic length of the RNA landscape scales linearly with  $n$ ,  $\lambda \approx \alpha n + \alpha_0$ , we find asymptotically

$$k = \frac{1}{\alpha} - 1 + O(n^{-1}). \quad (35)$$

Table VII compares the statistics of local optima in the

TABLE VII. Hypothetical  $k$  values and a comparison between the  $n$ - $k$  model and computed RNA data. Given are numbers from the older data set [22] and the newer one [23].  $k_{\text{pred}}$  is the predicted value of  $k$  for the RNA landscape as obtained from Eq. (34);  $f_{\text{pred}}$  is the predicted average scaled energy of a local optimum according to Eq. (28).  $f_{\text{LO}}$ ,  $f_{\text{LO}}^*$ , and  $f_{\text{LO}}^{**}$  are the numerical average scaled energies for random local optima and for local optima obtained from adaptive and gradient walks, respectively.

$n$	$k_{\text{pred}}$		$f_{\text{LO}}$		$f_{\text{LO}}^*$		$f_{\text{LO}}^{**}$		$f_{\text{pred}} = \eta$	
	[22]	[23]	[22]	[23]	[22]	[23]	[22]	[23]	[22]	[23]
GC										
20	6.58	6.99	1.41	1.85	1.99	2.57	2.15	2.69	0.73	0.72
25	6.74		1.66		2.57		2.68		0.73	
30	7.29	7.64	1.76	1.99	3.13	3.42	3.21	3.49	0.71	0.71
35	7.55		1.62		3.52		3.58		0.71	
40	7.89	7.85	1.68	2.50	3.75	3.88	3.75	3.96	0.70	0.70
50	8.64	8.94	2.14	2.91	4.17	4.41	4.10	4.48	0.69	0.70
60	8.73	7.98	2.04		4.53	4.88	4.48	4.92	0.68	0.68
GCAU										
30	3.14	3.55	1.95		5.62	6.69	6.27	7.38	1.10	1.07
40	2.86	3.62			6.53	7.83		8.54	1.13	1.07
50	2.51	3.08			7.17	8.54	8.02	9.39	1.16	1.11
60	2.17	3.38			7.83	9.37		10.13	1.19	1.08
70	2.29	2.29			8.73	10.18	9.51	10.97	1.18	1.12
80	2.29				9.19				1.18	

$n$ - $k$  landscape with the statistics of local optima in the RNA free-energy landscape. The ensembles of local optima in the RNA case have been obtained with gradient and adaptive walks, as detailed in Sec. IV E. The comparisons use the scaled quantities

$$f_{\text{LO}} = \frac{\langle F \rangle - \langle F_{\text{LO}} \rangle}{\sigma}. \quad (36)$$

Table VII shows that average free energies of local optima for the RNA landscapes, sampled with both gradient ( $f_{\text{LO}}^{**}$ ) and adaptive walks ( $f_{\text{LO}}^*$ ), differ significantly from those predicted by the  $n$ - $k$  model. Due to the sampling method, the local optima in the RNA case are obviously not a true random sample, but are biased towards those optima that have large basins of attractions. These optima usually have more extreme values than purely randomly sampled local optima. The comparison, therefore, is only limited. On the other hand, the  $n$ - $k$  model predicts a distribution of average energies of local optima that sharpens like  $\sigma_{\text{LO}}^2/\sigma^2 \sim 1/n$ , Eq. (29). If this were the case for the RNA landscape, we should observe no difference between our biased sample and a random sample. This is clearly not the case. In fact, the variances  $\sigma_{\text{LO}}^2$  (Table VII), scale linearly with chain length. In the GC case, we could accumulate a representative random sample of local optima (referred to as random optima in Sec. IV E). This allows a direct comparison with the prediction from the model as seen from Table VII;  $f_{\text{LO}}$  is still in disagreement with  $f_{\text{pred}}$ .

We find that some scaling properties are in common; in particular, the linear dependence of the walk length on  $n$  and the exponential decrease of the probability for

finding a local optimum  $p_{\text{LO}}$  as a function of  $n$ . Estimates for the latter can already be derived by assuming that in highly correlated landscapes there should be  $O(1)$  local optima in a path of radius  $\lambda$  [8], which, after some calculations, leads to

$$\ln p_{\text{LO}} = c_0 - n [\ln 2 + \ln(\kappa - 1)\alpha - 2(\alpha - 1/2)^2]. \quad (37)$$

By the same reasoning, the length of a gradient walk should be roughly

TABLE VIII. Scaling of the number of local optima and walk length for GC and AUGC landscapes. Shown are the slopes of linear fits to the calculated data (Table VII).  $n$ - $k$  model refers to the prediction from Sec. V. Also shown are the rough estimates obtained from Eqs. (37) and (38). Given are values of the older parameter set from Salser [22] and the newer set [23].

	$\ln p_{\text{LO}}/n$		$L_{\text{adap}}/n$		$L_{\text{grad}}/n$	
	[23]	[22]	[23]	[22]	[23]	[22]
GC						
GC landscape	-0.5	-0.12	0.15	0.18	0.10	0.12
$n$ - $k$ model	-0.25	-0.25	0.25	0.25	0.18	0.18
Eqs. (37), (38)	-0.35	-0.34			0.09	0.08
GCAU						
AUGC landscape	< -0.3	$\approx$ -0.2	0.44	0.32	0.26	0.20
$n$ - $k$ model	-0.60	-0.67	0.60	0.67	0.32	0.36
Eqs. (37), (38)	-0.87	-0.94			0.26	0.30

$$L_{\text{grad}} \approx \lambda \approx \alpha n. \quad (38)$$

Table VIII lists the scaled lengths of adaptive and gradient walks, as well as the logarithm of the probability of randomly hitting an optimum. The measured quantities are compared with the predictions from the  $n$ - $k$  model and the crude estimates based on Eqs. (37) and (38). The latter agree fairly well, in the case of gradient walks, with the observed data. The  $n$ - $k$  model is, at best, in the ballpark. The ratios  $L_{\text{adap}}/L_{\text{grad}}$ , Eq. (32), are pretty close to the predictions. For the GC alphabet we find 1.45, and for GCAU we find 1.57, versus a prediction of 1.45 and 1.85, respectively.

## VII. CONCLUSIONS

Among the most important steps in understanding evolutionary adaptation is the construction of a model landscape based on the proper abstractions of the adapting entities. In this paper, we explored in detail the statistics of a realistic and biologically motivated landscape induced by RNA folding. We compare its features with a widely used simple statistical model for rugged landscapes, known as the  $n$ - $k$  model.

RNA energy, as well as structure landscapes, were explored by numerically computing landscape autocorrelation functions, random-walk autocorrelation functions, and density surfaces as a function of the nucleotide alphabet and the sequence length. The main results can be summarized as follows.

### *RNA landscapes*

(1) The energy, as well as the structure autocorrelation function, is characterized to a reasonable approximation by one length scale.

(2) The characteristic length of the energy and of the structure landscape scale linearly with sequence length.

(3) The characteristic length for structures strongly depends on the nucleotide alphabet as follows.

(i) Binary sequences, AU or GC, have very short correlation lengths, indicating that they are very likely to change their structure with few changes in the underlying sequences.

(ii) GCXK sequences, with XK denoting two artificial nucleotides with the same pairing strength as GC, are less sensitive to changes than binary sequences.

(iii) Natural AUGC sequences are even less sensitive than GCXK. We have checked the influence of the non-Watson-Crick pair GU. Disabling GU pairs in AUGC sequences strongly influenced the energy autocorrelation (shorter correlation length), but had no or little effect on the structure autocorrelation. We conclude that the sensitivity difference between AUGC and GCXK is due to the unequal base stacking and pairing energies associated with GC and AU pairs.

(iv) ABCDEF sequences (GC pairing strength) have a very low sensitivity.

This suggests that a natural four-letter GCAU alpha-

bet is a good compromise between (a) enough structural variety to support biological function, and (b) sufficient, but not excessive, stability towards changes in the sequence.

(4) Exploration of the structure density surface suggests that there is a small region (compared to the diameter of the sequence space) around any random sequence, such that the sequences within that region fold into almost all minimum-free-energy structures.

### *Comparison with the $n$ - $k$ model*

The  $n$ - $k$  model is a powerful and flexible, yet simple, tool for generating scalar landscapes with prescribed correlation structure. Although both RNA and  $n$ - $k$  landscapes share some simple scaling behavior, they do not agree in important details concerning the statistics of local optima and the length of adaptive and gradient walks. We trace the disagreement between the fine structure of the two landscapes back to one basic difference. As detailed in Sec. IV E, the distribution of energy increments upon changes in one position are not properly described by a Gaussian in both binary alphabets and the natural alphabet. This is mainly due to a very high degree of neutrality, that is, neighboring sequences with identical minimum free energy or identical structure. With respect to energies, there are only a few thousand different values [30]. In the  $n$ - $k$  model, all values are pairwise distinct with probability one. Even a discretization of the  $n$ - $k$  model (e.g., cutting off all decimal places) would still remain Gaussian, without yielding a neutral-neighbor peak of the kind observed in RNA folding. The physical process of polynucleotide folding—as far as it is properly abstracted by the presently used algorithm—is not in the class defined by the  $n$ - $k$  model. The neutrality issue has profound effects on the number and the distribution of local optima as well as biased walks on both landscapes. These are the features in which the disagreement is most apparent. At the same time, these are also the features that are the most relevant to evolutionary optimization.

The autocorrelation function, Eq. (1), of the  $n$ - $k$  model is a single decaying exponential along a random walk. This is no longer the case for the landscape autocorrelation function, Eq. (2), as can be seen from Eqs. (23) and (24). However, the numerically computed (landscape) autocorrelation function of the RNA free-energy landscape is, to a good approximation, a single decaying exponential. For small distances [up to the “correlation length” (25)], we approximate Eq. (23) by an exponential. This is the basis for extracting a  $k \approx 7-8$  as the number of context sites influencing the energetic contribution of each position, independently of sequence length. This coincides roughly with the typical size of secondary structure elements [5], and one may speculate about the nature of this coincidence.

## ACKNOWLEDGMENTS

This work was partly supported by the Austrian Fonds zur Förderung der Wissenschaftlichen Forschung, Project Nos. S5305-PHY and P8526-MOB. The Sante Fe In-

stitute is sponsored by core funding from the U.S. Department of Energy (ER-FG05-88ER25054) and the National Science Foundation (PHY-8714918). Pedro Tarazona acknowledges support from the Direccion Central de Investigacion Cientifica y Tecnica, Project No. PB91-090.

**APPENDIX: DERIVATION  
OF THE AUTOCORRELATION FUNCTION  
FOR THE  $n$ - $k$  MODEL**

In case of the PR version, the chance that a site energy remains unchanged when a single digit is flipped equals  $1 - (k + 1)/n$ , independent of whether other bits have been flipped before. Therefore,

$$1 - h_{\text{PR}}(d) = \left[ 1 - \frac{k + 1}{n} \right]^d. \quad (\text{A1})$$

The RN version is only slightly more involved. A site energy  $f_i$  is unchanged if it is not one of the  $d$  positions that have been changed previously, and if it is not one of the  $k$  neighbors of any of the flipped sites. These events

are statistically independent, and thus

$$1 - h_{\text{RN}}(d) = \left[ 1 - \frac{d}{n} \right] \left[ 1 - \frac{k}{n - 1} \right]^d. \quad (\text{A2})$$

For the AN model, however, the situation is more complicated. First note that  $h(1) = (k + 1)/n$ , since, for a single bit flip, it does not matter whether the assignment of neighbors is random or not, because the contributions to the site energies are pairwise independent. Now suppose that the  $d \geq 2$  changes are labeled  $s_i$  such that

$$s_1 = 1 < s_2 < s_3 < \dots < s_d. \quad (\text{A3})$$

Evidently there are  $\binom{n-1}{d-1}$  such flipping patterns. Let us now calculate the probability that a particular pair of flips occurred  $l$  sites apart, i.e.,  $s_{i+1} - s_i = l$ . Suppose  $s_1$  and  $s_2$  are chosen this way. The remaining  $d - 2$  flips have therefore to be arranged within the remaining  $n - l - 1$  sites; there are  $\binom{n-l-1}{d-2}$  such arrangements. Therefore we obtain the probability that two flips are separated by  $l$  digits along the chain by

$$\phi_d(l) = \begin{cases} \binom{n-l-1}{d-2} / \binom{n-1}{d-1}, & 1 < l < n - 1, \quad 2 \leq d \leq n - l + 1, \\ 0, & \text{otherwise.} \end{cases}$$

If  $l \leq k$ , then  $l$  site energies are changed as a result of changing the corresponding position. Otherwise (all  $k + 1$  sites are changed). Thus the probability that a given site energy is changed by moving to Hamming distance  $d$  is

$$h_{\text{AN}}(d) = \frac{d}{n} \left[ \sum_{l=1}^k \phi_d(l)l + (k + 1) \left[ 1 - \sum_{l=1}^k \phi_d(l) \right] \right].$$

- \*Author to whom correspondence should be addressed, at Institut für Theoretische Chemie, Universität Wien Währingerstraße 17, A-1090 Vienna, Austria.
- [1] S. A. Kauffman and S. Levin, *J. Theor. Biol.* **128**, 11 (1987).
  - [2] M. Eigen, J. McCaskill, and P. Schuster, *Adv. Chem. Phys.* **75**, 149 (1989).
  - [3] E. D. Weinberger, *Biol. Cybern.* **63** (3), 25 (1990).
  - [4] W. Fontana, T. Griesmacher, W. Schnabl, P. F. Stadler, and P. Schuster, *Mh. Chem.* **122**, 795 (1991).
  - [5] W. Fontana, D. A. M. Konings, P. F. Stadler, and P. Schuster, *Biopolymers* (to be published).
  - [6] E. D. Weinberger and P. F. Stadler (to be published).
  - [7] B. Derrida, *Phys. Rev. B* **24**, 2613 (1981).
  - [8] P. F. Stadler and W. Schnabl, *Phys. Lett. A* **161**, 337 (1992).
  - [9] P. F. Stadler (unpublished).
  - [10] P. F. Stadler and R. Happel, *J. Phys. A* **25**, 3103 (1992).
  - [11] E. D. Weinberger (unpublished).
  - [12] E. L. Lawler, J. K. Lenstra, A. H. G. Rinnoy Kan, and D. B. Shmoys, *The Traveling Salesman Problem. A Guided Tour of Combinatorial Optimization* (Wiley, New York, 1985).
  - [13] M. Mézard and G. Parisi, *Europhys. Lett.* **2**, 913 (1986).
  - [14] Y. Fu and P. W. Anderson, *J. Phys. A* **19**, 1605 (1986).
  - [15] J. Bernasconi, *J. Phys. (Paris)* **48**, 559 (1987).

- [16] D. Sherrington and S. Kirkpatrick, *Phys. Rev. Lett.* **35**, 1792 (1975).
- [17] C. Amitrano, L. Peliti, and M. Saber, *J. Mol. Evol.* **29**, 513 (1989).
- [18] R. Nussinov, G. Pieczenik, J. R. Griggs, and D. J. Kleitman, *SIAM J. Appl. Math.* **35**, 68 (1978).
- [19] M. S. Waterman and T. F. Smith, *Math. Biosci.* **42**, 257 (1978).
- [20] M. Zuker and P. Stiegler, *Nucleic Acids Res.* **9**, 133 (1981).
- [21] M. Zuker and D. Sankoff, *Bull. Math. Biol.* **46**, 591 (1984).
- [22] W. Salsler, *Cold Spring Harbor Symp. Quant. Biol.* **42**, 985 (1977).
- [23] S. M. Freier, R. Kierzek, J. A. Jaeger, N. Sugimoto, M. H. Caruthers, T. Neilson, and D. H. Turner, *Biochemistry* **83**, 9373 (1986).
- [24] J. A. Jaeger, D. H. Turner, and M. Zuker, *Proc. Natl. Acad. Sci. U.S.A.* **86**, 7706 (1989).
- [25] I. Hofacker, P. Schuster, and P. F. Stadler (unpublished).
- [26] D. A. M. Konings and P. Hogeweg, *J. Mol. Biol.* **207**, 597 (1989).
- [27] P. Hogeweg and P. Hesper, *Nucleic Acids Res.* **12**, 67 (1984).
- [28] B. A. Shapiro, *CABIOS* **4**, 381 (1988).
- [29] A. S. Perelson and G. F. Oster, *J. Theor. Biol.* **81**, 645 (1979).



- [30] W. Fontana, W. Schnabl, and P. Schuster, *Phys. Rev. A* **40**, 3301 (1989).
- [31] S. A. Kauffman and E. D. Weinberger, *J. Theor. Biol.* **141**, 211 (1989).
- [32] S. A. Kauffman, E. D. Weinberger, and A. S. Perelson, in *Theoretical Immunology, Part I Santa Fe Institute Studies in the Sciences of Complexity*, edited by A. S. Perelson (Addison-Wesley, Reading, MA, 1988).
- [33] E. D. Weinberger, *J. Theor. Biol.* **134**, 125 (1988).
- [34] C. A. Macken and A. S. Perelson, *Proc. Natl. Acad. Sci. U.S.A.* **86**, 6191 (1989).
- [35] C. A. Macken, P. S. Hagan, and A. S. Perelson, *SIAM J. Appl. Math.* **51**, 799 (1991).
- [36] E. D. Weinberger, *Phys. Rev. A* **44**, 6399 (1991).